# Machine Learning in the audio domain
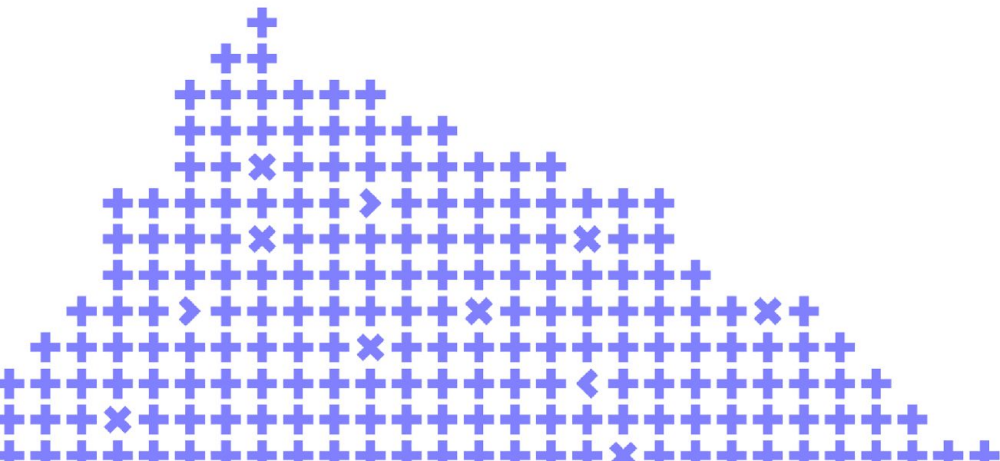
**When the neural network is overkill or where are the limits of lightweight models**
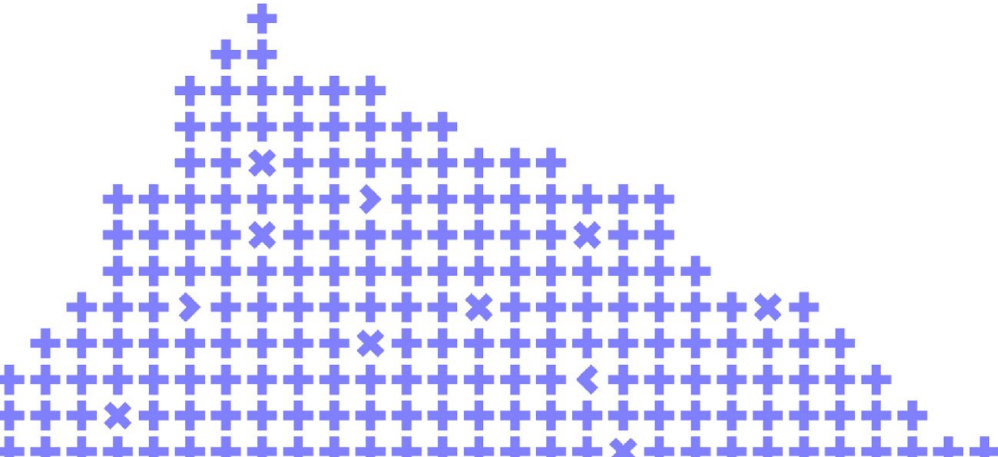
# Roman Smirnov

## Machine Learning Engineer

- **1 year at Exness**
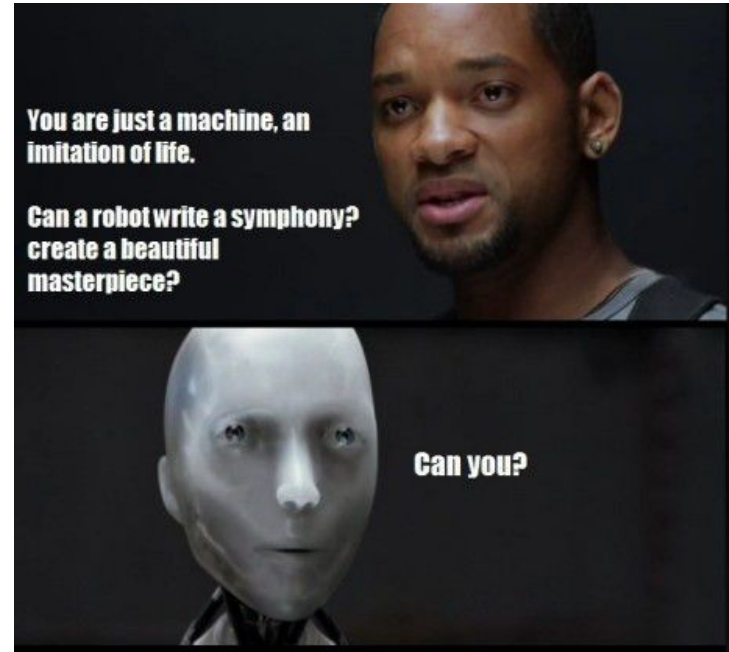- **7 years as an MLE/Tech Lead at Skyeng & MSU Labs**

# Table of contents

- Audio domain
- Problem statement
- Audio-domain tasks
- Experiments & Business
- Training process
- Results
- Business outcomes

# **Audio domain**

Analyzing sounds using
neural networks

- Call Centers
- Virtual Assistants
- Speech and music
  generation

# Problem Statement

# Data modality and SotA

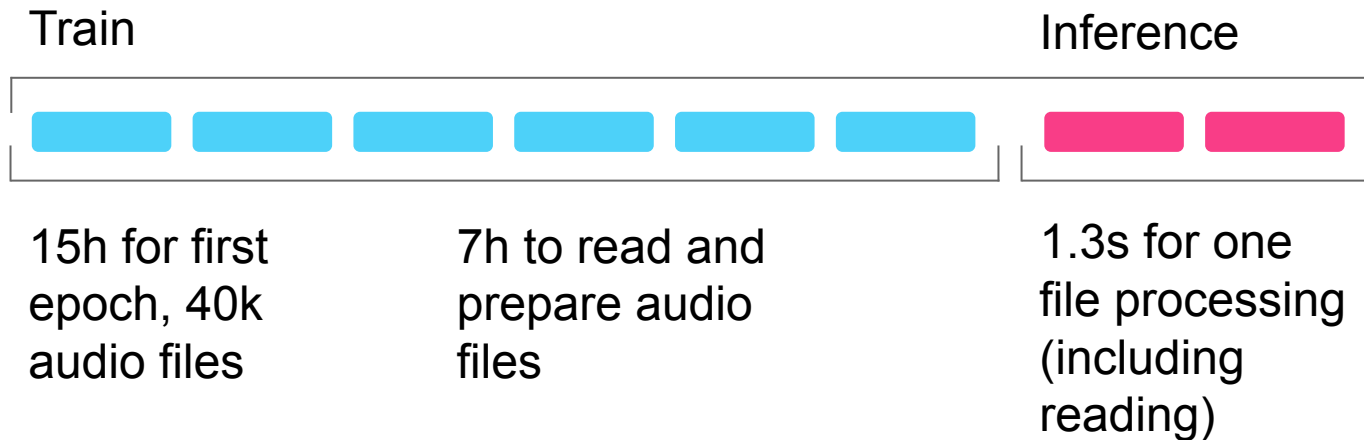When working with media data, we usually use large neural networks.

But:

- They are resource-heavy
- They are slower than light models

Usually it is **Transformers models**

Pic source https://comicvine.gamespot.com/

# Time for SotA

It takes **LONG** time to train SotA model for the task we'll discuss today

Train

Inference

15h for first epoch, 40k audio files

7h to read and prepare audio files

1.3s for one file processing (including reading)

GPU: RTX 3090Ti, 24GB

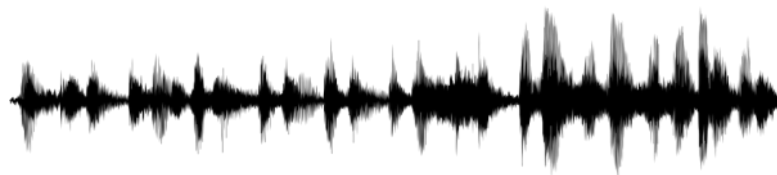# Audio-domain tasks
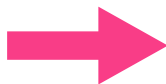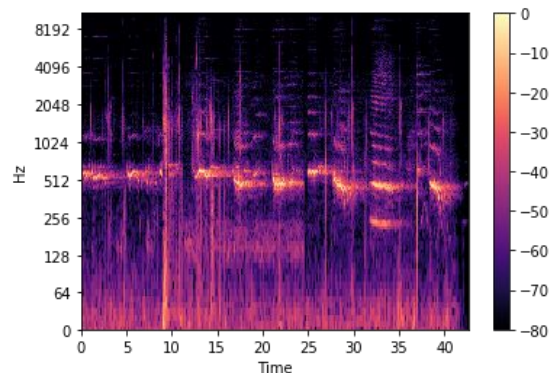
Understanding and Generation

# Understanding (Classification)

- Classification, e.g. emotions classification
- "Token classification"—voice activity detection
- Speaker separation, user verification
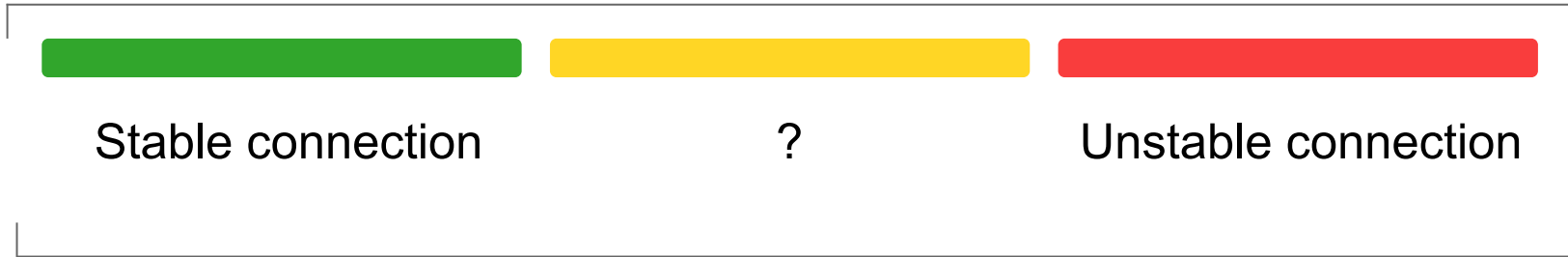- Automatic Speech Recognition (ASR)

# Generation



- Voice cloning
- Text to speech
- Speech and music generation
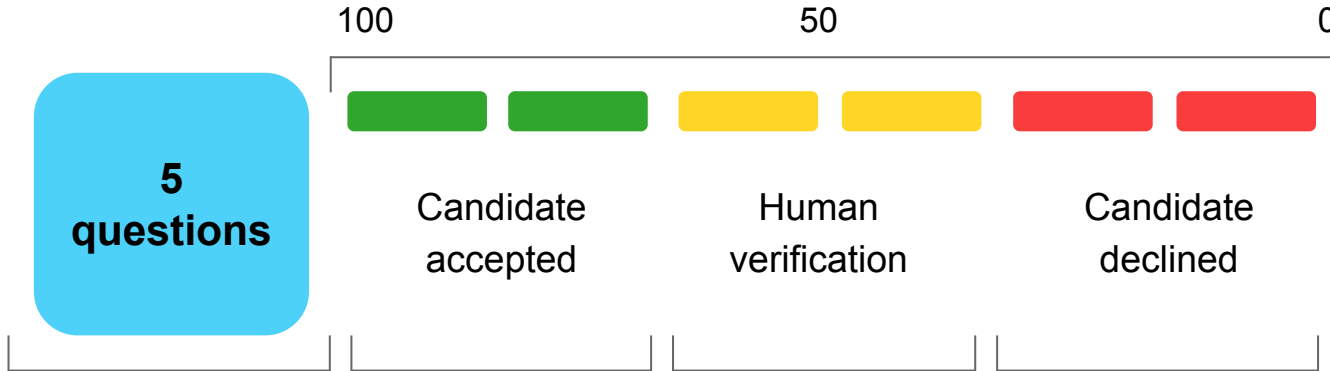
# Experiments & Business

# Classification

Evaluation of the call quality for the call center:



Stable connection    ?    Unstable connection

# Regression

Interviewing a candidate for an English Teacher position

| 100 | 50 | 0 |
|---|---|---|

**5 questions**

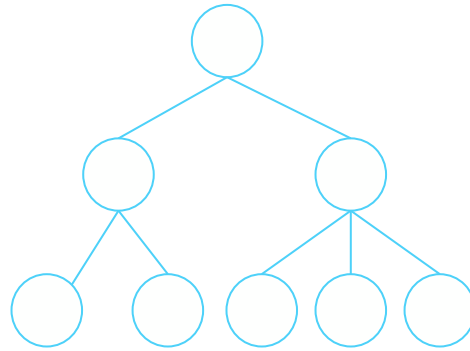Candidate accepted

Human verification

Candidate declined

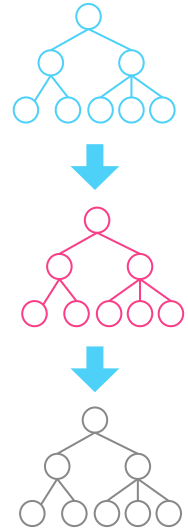Fluency and pronunciation are evaluated separately using 5-points-scale

# Gradient Boosting on Decision Trees (GB on DT)

- Fast
- Work on statistics aggregates of audio
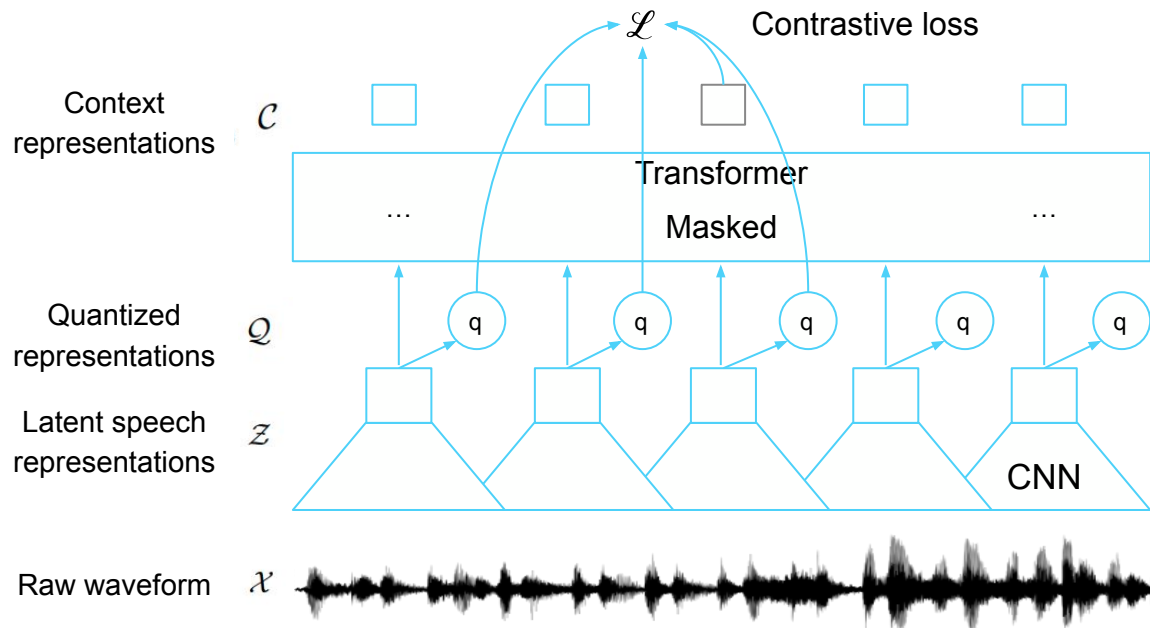- But is it accurate?

Single Decision Tree

Gradient Boosted Trees

# Wav2Vec2

- SotA for speech recognition
- Transformer model
- CNNs allows to encode local features
- But is it fast?



Context representations $\mathcal{C}$

Contrastive loss $\mathcal{L}$

Transformer

Masked

Quantized representations $\mathcal{Q}$

Latent speech representations $\mathcal{Z}$

CNN

Raw waveform $\mathcal{X}$

# Training Process

# Preprocessing

**For GB on DT**

- Collect amplitude and melspectrograms statistics
- Mean, median, min, max, std, kurtosis, skew… over mel for every frequency and raw wav
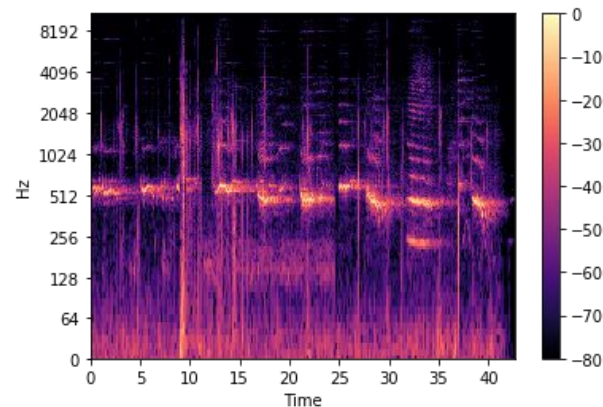
**For Wav2Vec2**

- Resample to 16kHz
- Just truncation

17

# What is melspectrogram?

| Raw audio | Melspectrogram |
|---|---|

# Training

| For GB on DT | For Wav2Vec2 |
| --- | --- |
| • Catboost<br><br>• Small trees (depth 6)<br><br>• 100-10000 iterations | • Wav2Vec2 base<br><br>• Freeze feature encoder<br><br>• Learning rate schedulers<br><br>• 10 epochs |

# Time metrics. Classification

## For GB on DT

- Train: 45 min to read data
- 0.5s to train on GPU
- 5k audio-files
- Inference: 5 files, 2s to read data, 0.1s to inference on CPU

## For Wav2Vec2

- Train: 45 min to read data
- 5h to train on GPU (10 epochs)
- 5k audio-files, batch size 6
- Inference: 5 files, 2s to read data, 105s to inference on CPU

# Time metrics. Regression

## For GB on DT

- Train: 7h to read data
- 5 min to train on GPU
- 40k audio-files
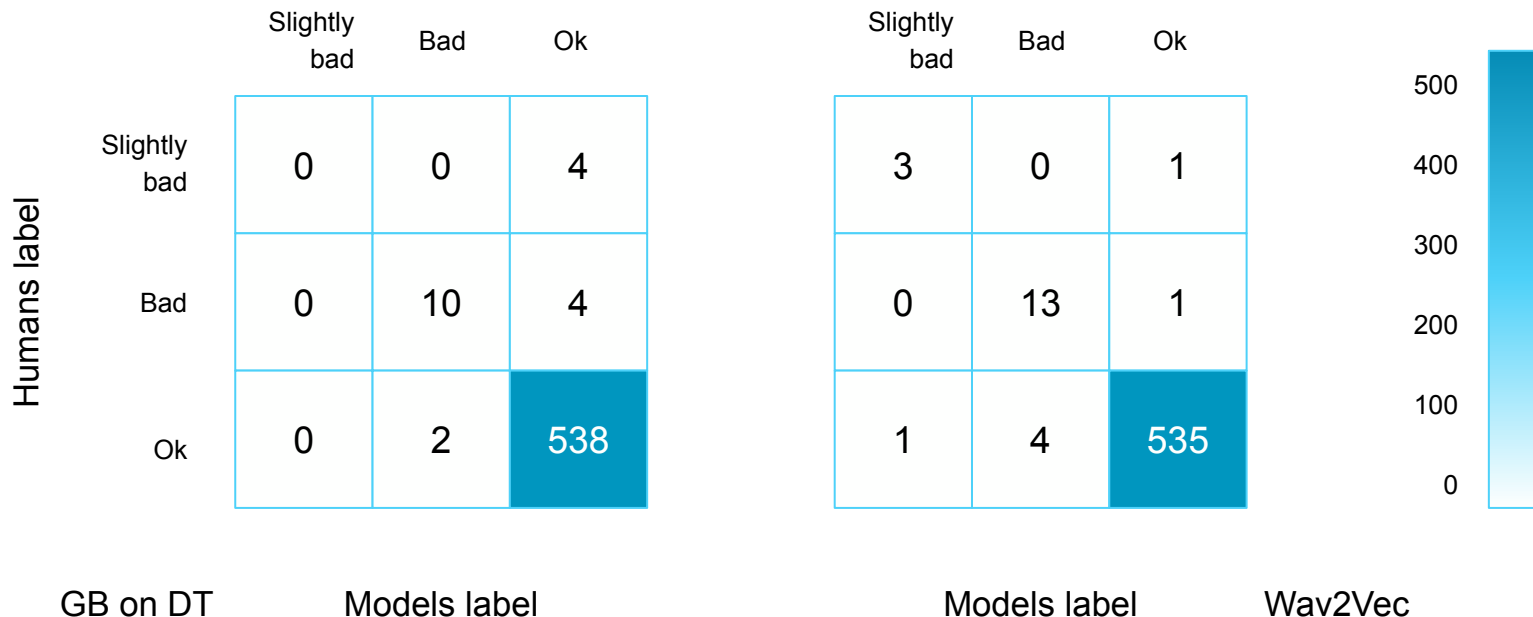- Inference: 5 files, 2s to read data, 0.1s to inference on CPU

## For Wav2Vec2

- Train: 7h to read data
- 90h to train on GPU (10 epochs)
- 40k audio-files, batch size 6
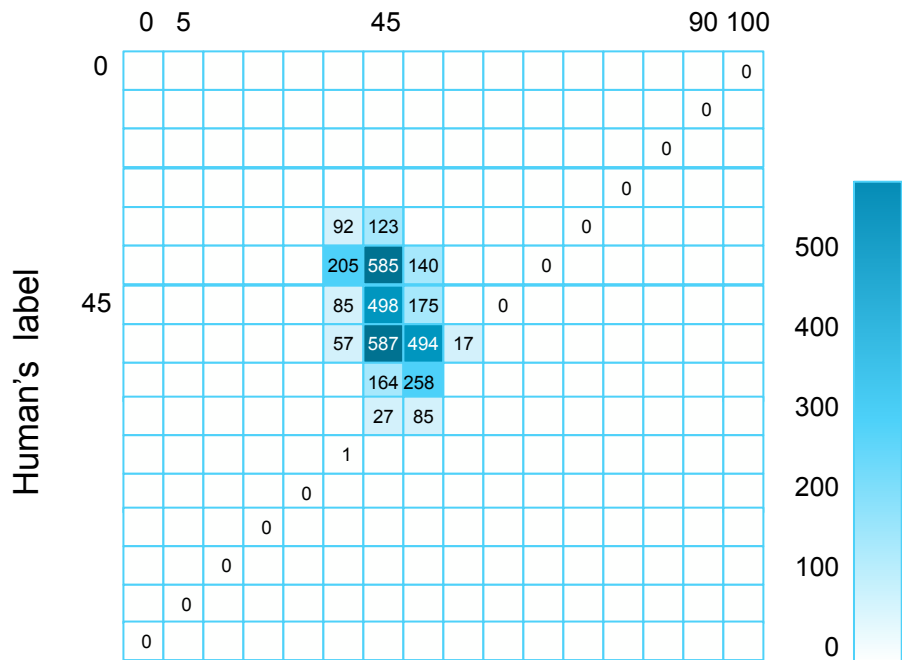- Inference: 5 files, 2s to read data, 105s to inference on CPU

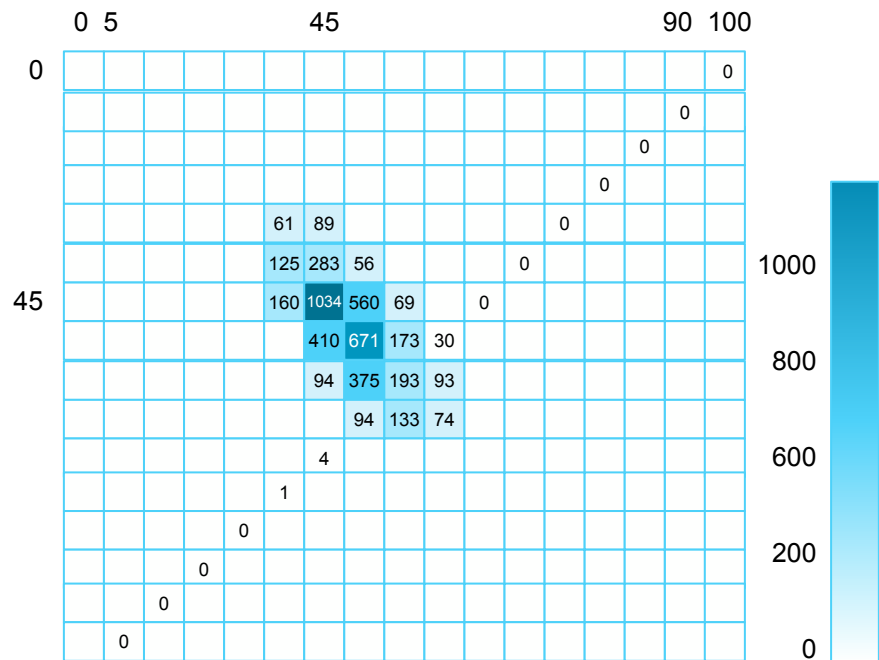# Results

# Results. Classification task

**Natural imbalance!**



GB on DT      Models label

Wav2Vec      Models label

# Results. Regression task

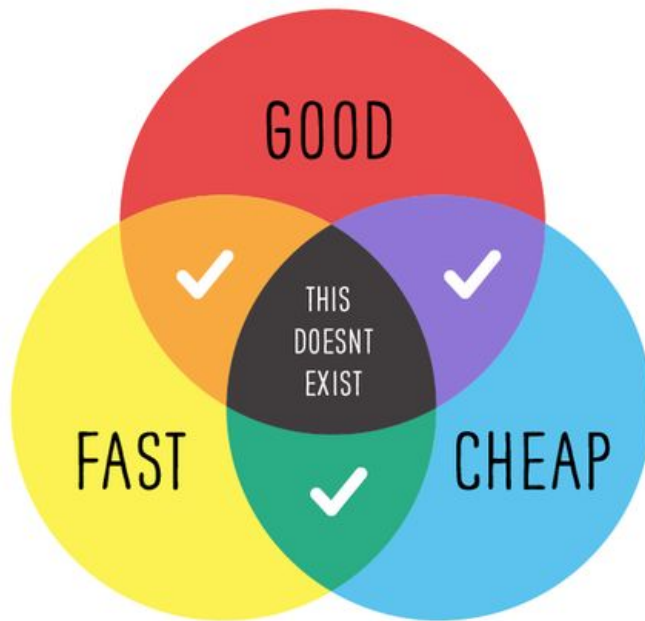

GB on DT — Model's label

Wav2Vec — Model's label

# Business Outcomes

# Outcomes

**Time and $ / accuracy trade-off**

- We took GB on DT for classification
- We took Wav2Vec2 for regression

**Where are the limits and what is overkill?**



GOOD

FAST

CHEAP

THIS DOESNT EXIST

HEINLEY

# Thank you for your attention!

# Vote for my talk

Roman Smirnov

rgsmirnov

Exness Tech Blog